

This Page Is Inserted by IFW Operations  
and is not a part of the Official Record

## **BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning documents *will not* correct images,  
please do not report the images to the  
Image Problem Mailbox.**

**THIS PAGE BLANK (USPTO)**

4 日本国特許庁

PATENT OFFICE  
JAPANESE GOVERNMENT

02.06.00

JP00/3623

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出願年月日

Date of Application:

1999年 6月17日

27 JUL 2000

出願番号

Application Number:

平成11年特許願第171723号

出願人

Applicant(s):

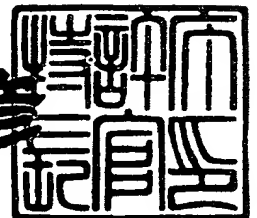
セイコーエプソン株式会社

PRIORITY  
DOCUMENTSUBMITTED OR TRANSMITTED IN  
COMPLIANCE WITH RULE 17.1(a) OR (b)

2000年 6月29日

特許庁長官  
Commissioner,  
Patent Office

近藤 隆彦



出証番号 出証特2000-3052043

【書類名】 特許願

【整理番号】 J0074154

【提出日】 平成11年 6月17日

【あて先】 特許庁長官 殿

【国際特許分類】 G06F 7/24

【発明の名称】 情報分類方法及び情報分類装置並びに情報分類処理プログラムを記録した記録媒体

【請求項の数】 21

【発明者】

    【住所又は居所】 長野県諏訪市大和3丁目3番5号 セイコーエプソン株式会社内

    【氏名】 長石 道博

【特許出願人】

    【識別番号】 000002369

    【氏名又は名称】 セイコーエプソン株式会社

    【代表者】 安川 英昭

【代理人】

    【識別番号】 100093388

    【弁理士】

    【氏名又は名称】 鈴木 喜三郎

    【連絡先】 0 2 6 6 - 5 2 - 3 1 3 9

【選任した代理人】

    【識別番号】 100095728

    【弁理士】

    【氏名又は名称】 上柳 雅誉

【選任した代理人】

    【識別番号】 100107261

    【弁理士】

    【氏名又は名称】 須澤 修

【手数料の表示】

【予納台帳番号】 013044

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9711684

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 情報分類方法及び情報分類装置並びに情報分類処理プログラム  
を記録した記録媒体

【特許請求の範囲】

【請求項 1】 検索サービスにユーザからの検索要求が与えられることによって検索された複数の文書に対し、それぞれの文書の共通性に基づいてクラスタリング処理し、それによって得られたクラスタリング結果に対し、検索されたそれぞれの文書対応に付された検索要求との適合性を示す値（スコアという）を用いて、前記クラスタリング処理によって得られたそれぞれのクラスタの順位を再構成し、そのクラスタ順位が再構成されたクラスタリング結果を出力することを特徴とする情報分類方法。

【請求項 2】 前記それぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書対応に付されたスコアの平均値をそれぞれのクラスタごとに求め、クラスタごとの平均値をそれぞれのクラスタのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成することを特徴とする請求項 1 記載の情報分類方法。

【請求項 3】 前記それぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書対応に付されたスコアの最大値をそれぞれのクラスタごとに得て、そのクラスタごとのスコアの最大値をそれぞれのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成することを特徴とする請求項 1 記載の情報分類方法。

【請求項 4】 前記それぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書がそれぞれの文書対応に付されたスコアの大きい順に並べられている場合、その中央または中央付近に位置するスコアをそれぞれのクラスタごとに得て、そのクラスタごとの中央または中央付近に位置するスコアをそれぞれのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成することを特徴とする請求項 1 記載の情報分類方法。

【請求項 5】 前記クラスタリング処理を複数の検索サービスによって得られた検索結果に対応して行うことを可能とする場合、前記クラスタの順位を再構

成するためのクラスタスコアを求める処理は、複数の検索サービスそれぞれに対応して行うことを特徴とする請求項 2 から 4 のいずれか 1 項に記載の情報分類方法。

【請求項 6】 前記クラスタリング処理は、それぞれの文書のタイトルを検出し、そのタイトルに含まれる特徴的な用語を特徴要素として抽出し、抽出された特徴要素に基づいて行うことを特徴とする請求項 1 から 5 のいずれか 1 項に記載の情報分類方法。

【請求項 7】 前記クラスタ順位が再構成されたクラスタリング結果の出力の仕方は、クラスタスコアの高いクラスタ順に表示し、クラスタスコアが同じであるクラスタが存在する場合には、クラスタ内の文書数の多いクラスタを高順位とすることを特徴とする請求項 1 から 6 のいずれか 1 項に記載の情報分類方法。

【請求項 8】 検索サービスにユーザからの検索要求が与えられることによって検索された複数の文書に対し、それぞれの文書の共通性に基づいてクラスタリング処理するクラスタリングモジュールと、

このクラスタリングモジュールによって得られたクラスタリング結果に対し、検索されたそれぞれの文書対応に付された検索要求との適合性を示す値（スコアという）を用いて、前記クラスタリング処理によって得られたそれぞれのクラスタの順位を再構成し、そのクラスタ順位が再構成されたクラスタリング結果を出力するクラスタ順位設定モジュールと、

を有することを特徴とする情報分類装置。

【請求項 9】 前記クラスタ順位設定モジュールが行うそれぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書対応に付されたスコアの平均値をそれぞれのクラスタごとに求め、クラスタごとの平均値をそれぞれのクラスタのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成することを特徴とする請求項 8 記載の情報分類装置。

【請求項 10】 前記クラスタ順位設定モジュールが行うそれぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書対応に付されたスコアの最大値をそれぞれのクラスタごとに得て、そのクラスタごとのスコアの最大値をそれぞれのクラスタスコアとし、そのクラスタスコアによ

って、クラスタの順位を再構成することを特徴とする請求項 8 記載の情報分類装置。

【請求項 11】 前記クラスタ順位設定モジュールが行うそれぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書がそれぞれの文書対応に付されたスコアの大きい順に並べられている場合、その中央または中央付近に位置するスコアをそれぞれのクラスタごとに得て、そのクラスタごとの中央または中央付近に位置するスコアをそれぞれのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成することを特徴とする請求項 8 記載の情報分類装置。

【請求項 12】 前記クラスタリング処理を複数の検索サービスによって得られた検索結果に対応して行うことを可能とする場合、前記クラスタの順位を再構成するためのクラスタスコアを求める処理は、複数の検索サービスそれぞれに対応して行うことを特徴とする請求項 9 から 11 のいずれか 1 項に記載の情報分類装置。

【請求項 13】 前記クラスタリングモジュールが行うクラスタリング処理は、それぞれの文書のタイトルを検出し、そのタイトルに含まれる特徴的な用語を特徴要素として抽出し、抽出された特徴要素に基づいて行うことを特徴とする請求項 8 から 12 のいずれか 1 項に記載の情報分類装置。

【請求項 14】 前記クラスタ順位が再構成されたクラスタリング結果の出力の仕方は、クラスタスコアの高いクラスタ順に表示し、クラスタスコアが同じであるクラスタが存在する場合には、クラスタ内の文書数の多いクラスタを高順位とすることを特徴とする請求項 8 から 13 のいずれか 1 項に記載の情報分類装置。

【請求項 15】 検索サービスにユーザからの検索要求が与えられることによって検索された複数の文書をクラスタリング処理してそのクラスタリング結果を出力する情報分類処理プログラムを記録した記録媒体であって、その情報分類処理プログラムは、

検索サービスによって検索された複数の文書に対し、それぞれの文書の共通性に基づいてクラスタリング処理する手順と、



これによって得られたクラスタリング結果に対し、検索されたそれぞれの文書対応に付された検索要求との適合性を示す値（スコアという）を用いて、前記クラスタリング処理によって得られたそれぞれのクラスタの順位を再構成し、そのクラスタ順位が再構成されたクラスタリング結果を出力する手順と、

を含むことを特徴とする情報分類処理プログラムを記録した記録媒体。

【請求項 16】 前記それぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書対応に付されたスコアの平均値をそれぞれのクラスタごとに求め、クラスタごとの平均値をそれぞれのクラスタのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成することを特徴とする請求項 15 記載の情報分類処理プログラムを記録した記録媒体。

【請求項 17】 前記それぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書対応に付されたスコアの最大値をそれぞれのクラスタごとに得て、そのクラスタごとのスコアの最大値をそれぞれのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成することを特徴とする請求項 15 記載の情報分類処理プログラムを記録した記録媒体。

【請求項 18】 前記それぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書がそれぞれの文書対応に付されたスコアの大きい順に並べられている場合、その中央または中央付近に位置するスコアをそれぞれのクラスタごとに得て、そのクラスタごとの中央または中央付近に位置するスコアをそれぞれのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成することを特徴とする請求項 15 記載の情報分類処理プログラムを記録した記録媒体。

【請求項 19】 前記クラスタリング処理を複数の検索サービスによって得られた検索結果に対応して行うことを可能とする場合、前記クラスタの順位を再構成するためのクラスタスコアを求める処理は、複数の検索サービスそれぞれに対応して行うことを特徴とする請求項 16 から 18 のいずれか 1 項に記載の情報分類処理プログラムを記録した記録媒体。

【請求項 2 0】 前記クラスタリング処理は、それぞれの文書のタイトルを検出し、そのタイトルに含まれる特徴的な用語を特徴要素として抽出し、抽出された特徴要素に基づいて行うことを特徴とする請求項 1 5 から 1 9 のいずれか 1 項に記載の情報分類処理プログラムを記録した記録媒体。

【請求項 2 1】 前記クラスタ順位が再構成されたクラスタリング結果の出力の仕方は、クラスタスコアの高いクラスタ順に表示し、クラスタスコアが同じであるクラスタが存在する場合には、クラスタ内の文書数の多いクラスタを高順位とすることを特徴とする請求項 1 5 から 2 0 のいずれか 1 項に記載の情報分類処理プログラムを記録した記録媒体。

【発明の詳細な説明】

【0 0 0 1】

【発明の属する技術分野】

本発明は汎用の検索サービスで検索された結果に対しクラスタリング処理を施し、かつ、クラスタリング処理によって得られたクラスタの順位を再構成して表示することで、ユーザの検索要求に沿ったクラスタリング結果を提示できるようにした情報分類方法及び情報分類装置並びに情報分類処理プログラムを記録した記録媒体に関する。

【0 0 0 2】

【従来の技術】

ネットワーク上に存在する膨大な量の情報の中からユーザの所望とする情報を検索する場合、検索サービスの存在は重要である。たとえば、インターネットで web ページを検索する際、ユーザは、幾つかの検索サービスの中から任意の検索サービスを選び、自分の欲しい情報を得るための検索要求としてのキーワードを入力する。これによって、検索サービス側では、入力されたキーワードに基づいて情報検索を行って、その検索結果をユーザに提示する。

【0 0 0 3】

しかし、検索サービスによって検索される情報は膨大な量となることも多く、その中からユーザの本当に欲しい情報を見つけるのは非常に大変である。近年、web ページは増大の一途をたどっているため、検索された多数の情報を如何にユ

ーザにわかりやすく提示するかが大きな課題となっている。

#### 【0004】

最近では、検索された情報をユーザが見やすい形に加工して提示する手法も開発され実用化されつつある。たとえば、ユーザの入力したキーワードで検索された結果から得られるキーワードを用いて再検索することで、情報の絞り込みを行い、ユーザの所望とするwebページを見つけやすくする方法がある。つまり、検索によって得られる検索結果の集合を特徴づけるキーワードを抽出して、ユーザの本当に欲しい情報の集合に収束させる方法である。

#### 【0005】

このように、膨大な情報の中から、似た情報の集合を見つけることをクラスタリングという。情報処理ではこのクラスタリングはよく知られた手法であり、膨大な文書を分類する場合などに広く使われている。

#### 【0006】

一般にクラスタリングは、文書の一部、たとえば、webページの場合、文書のタイトル（見出し）などに含まれる特徴的な用語を特徴要素として抽出し、その特徴要素を用いてクラスタリングすることが行われている。勿論、文書全体を解析して特徴要素を抽出し、抽出された特徴要素の基づいてクラスタリングすることも可能であるが、演算量が多く処理に時間がかかる問題がある。

#### 【0007】

特に、汎用の検索サービスを利用してwebページの検索を行ったのちに、その検索結果の後処理としてクラスタリングを行う場合は、リアルタイムでクラスタリング結果を出力させる必要があり、高速な処理が要求されるので、文書全体を解析してクラスタリングする手法は現実的ではない。その点、上述した文書のタイトルから特徴要素してクラスタリングする手法は、演算量も少なく処理時間も大幅に短縮され、しかも、文書のタイトルは文書の中身を適切に反映したものが多くことから、適切なクラスタリングが可能となるなどの特徴を有している。このため、今後もこの手法は広く使用されるものと考えられる。

## 【0008】

## 【発明が解決しようとする課題】

しかし、上述した文書のタイトルから特徴要素してクラスタリングする手法は、演算量や処理時間の面で優れ、適切なクラスタリングが可能となるとはいつても、クラスタリングを行うための情報量は文書全体から見れば少ないので、全てが適切にクラスタリングされるとは限らない。特に、タイトルが文書の内容を適切に表していなかったり、文書内容とは大きくかけ離れた奇抜なタイトルが付けられていたりする場合もある。このような場合には、クラスタリング精度は大きく低下し、良好なクラスタリング結果は得られないことになる。

## 【0009】

また、特徴要素を抽出して、その特徴要素に基づいてクラスタリングする手法は、特徴要素の頻度などを調べ、それによって、機械的に文書を分類してクラスタリングするのが一般的である。このようなクラスタリングでは、文書の意味を解析しているわけではないので、得られたクラスター（クラスタリングされることによって得られる文書群の1つの集合）が必ずしも意味的な共通性のある文書の集合となるとは限らない。

## 【0010】

このように、これまでのクラスタリングは、多数の文書をそのタイトルなどから特徴要素を抽出するなどして、ある程度の適切さを有したクラスタリングがなされるものの、そのクラスタリング結果は、ユーザ側から見て満足度の高いものとはならない場合も多い。

## 【0011】

たとえば、ユーザが「半導体」というキーワードをある検索サービスに与えたとき、その検索サービスから得られた多数の文書に対しクラスタリング処理した結果を例にして考える。このクラスタリング結果の一例を図2に示す。

## 【0012】

この図2で示されるクラスタリング結果は、クラスタリングによって得られた各クラスターの名称（クラスター名という）と、それぞれのクラスターに属する検索結果、それぞれのクラスターに幾つの文書が含まれているかを示す文書数などが一覧

表形式で示されている。

【0013】

なお、検索結果としては、検索された文書のタイトルとそのタイトルに付されたスコアと呼ばれる数値（たとえば、表の最上段にある「東北エプソンの概要」という文書のタイトルに付された「133」という数値）が表示される。このスコアは、与えられたキーワードに対しその検索サービスが独自の方法で検索し、検索された文書対応に付けられた値であり、一般には、与えられたキーワードとそれぞれの文書との適合度を示す客観的な尺度として用いられ、汎用の検索サービスによる検索結果にはこのスコアが付されるのが普通である。このスコアは、検索方式により計算の仕方や値の考え方などが異なるものの、一般に、値が大きいほど、与えられたキーワードに適合する内容を有する文書であるといえる。

【0014】

この図2で示される検索結果において、それぞれのクラスタは文書数の多い順に記述されている。すなわち、「概要」というクラスタ名を持つクラスタ（以下、概要クラスタという）は、そのクラスタ内の文書数が16個であり、この場合、最も多いので、トップに位置づけられており、続いて、「LP（レーザプリンタ）」というクラスタ名を持つクラスタ（以下、LPクラスタという）は、そのクラスタ内の文書数が16個で概要クラスタと同数であるが、ここでは2番目に記述されており、さらに続いて、「仕様」というクラスタ名を持つクラスタ（以下、仕様クラスタという）は、そのクラスタ内の文書数が14で3番目に記述されており、以下、「デバイス」というクラスタ名を持つクラスタ（以下、デバイスクラスタという）は、そのクラスタ内の文書数が9個で4番目、「半導体」というクラスタ名を持つクラスタ（以下、半導体クラスタという）は、そのクラスタ内の文書数が7個で5番目、「電子」というクラスタ名を持つクラスタ（以下、電子クラスタという）は、そのクラスタ内の文書数が4個で6番目に記述されている。

【0015】

この図2をみると、ユーザの入力したキーワード「半導体」に対して検索された膨大な文書をクラスタリングすると、概要クラスタ、LPクラスタ、仕様クラ

スタ、デバイスクラスタ、半導体クラスタ、電子クラスタなどのクラスタに分類されることがわかる。そして、それぞれのクラスタに含まれる文書数によって図2のようなクラスタ順位での表示がなされる。

#### 【0016】

ところで、検索結果をクラスタリングしてそのクラスタリング結果を図2のような一覧表形式でユーザ側端末に表示する場合は、ユーザの本当に必要とする情報を一覧表の上位に位置させる方がユーザにとっては見易いものとなるのは言うまでもない。

#### 【0017】

しかしながら、図2の例では、ユーザの欲する情報が最上位にきているとは思えない。すなわち、概要クラスタ、LPクラスタ、仕様クラスタに含まれる文書の見出しから判断すれば、それらのクラスタには、多種のレーザプリンタに関する仕様概要に関する文書が多く、半導体に直接関係している内容を有する文書は少ないと思われる。この場合、ユーザの与えたキーワードに対してそれに関連しそうな情報が多く含まれそうなクラスタとしては、デバイスクラスタ、半導体クラスタ、電子クラスタであって、これらのクラスタが上位に来ることが望ましいが、この検索結果のクラスタリング結果では、それらのクラスタは下位に位置した状態となっている。

#### 【0018】

この図2で示されるクラスタリング結果はほんの一例であるが、従来のクラスタリング処理は、膨大な検索結果がクラスタリングされて出力されるものの、そのクラスタリングによって得られるそれぞれのクラスタのクラスタ順位はユーザの与えたキーワードに対し適切な結果とは言えない内容であることがしばしば起こる。

#### 【0019】

そこで本発明は、汎用の検索サービスで得られた検索結果をクラスタリング処理を施し、かつ、クラスタリングによって得られたクラスタの順位を再構成して表示することで、ユーザの検索要求に沿ったクラスタリング結果を提示できるようにすることを目的としている。

【 0 0 2 0 】

## 【課題を解決するための手段】

前述の目的を達成するために、本発明の情報分類方法は、検索サービスにユーザからの検索要求が与えられることによって検索された複数の文書に対し、それぞれの文書の共通性に基づいてクラスタリング処理し、それによって得られたクラスタリング結果に対し、検索されたそれぞれの文書対応に付された検索要求との適合性を示す値（スコアという）を用いて、前記クラスタリング処理によって得られたそれぞれのクラスタの順位を再構成し、そのクラスタ順位が再構成されたクラスタリング結果を出力するようにしている。

【 0 0 2 1 】

また、本発明の情報分類装置は、検索サービスにユーザからの検索要求が与えられることによって検索された複数の文書に対し、それぞれの文書の共通性に基づいてクラスタリング処理するクラスタリングモジュールと、このクラスタリングモジュールによって得られたクラスタリング結果に対し、検索されたそれぞれの文書対応に付された検索要求との適合性を示す値（スコアという）を用いて、前記クラスタリング処理によって得られたそれぞれのクラスタの順位を再構成し、そのクラスタ順位が再構成されたクラスタリング結果を出力するクラスタ順位設定モジュールとを有する構成としている。

【 0 0 2 2 】

また、本発明の情報分類処理プログラムを記録した記録媒体は、検索サービスにユーザからの検索要求が与えられることによって検索された複数の文書をクラスタリング処理してそのクラスタリング結果を出力する情報分類処理プログラムを記録した記録媒体であって、その情報分類処理プログラムは、検索サービスによって検索された複数の文書に対し、それぞれの文書の共通性に基づいてクラスタリング処理する手順と、これによって得られたクラスタリング結果に対し、検索されたそれぞれの文書対応に付された検索要求との適合性を示す値（スコアという）を用いて、前記クラスタリング処理によって得られたそれぞれのクラスタの順位を再構成し、そのクラスタ順位が再構成されたクラスタリング結果を出力する手順とを含むようにしている。

## 【0023】

これら各発明において、前記それぞれのクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書対応に付されたスコアの平均値をそれぞれのクラスタごとに求め、クラスタごとの平均値をそれぞれのクラスタのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成するようにしている。

## 【0024】

このクラスタの順位を再構成する処理は、それぞれのクラスタに含まれるそれぞれの文書対応に付されたスコアの最大値をそれぞれのクラスタごとに得て、そのクラスタごとのスコアの最大値をそれぞれのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成するようにしてもよく、それぞれのクラスタに含まれるそれぞれの文書がそれぞれの文書対応に付されたスコアの大きい順に並べられている場合、その中央または中央付近に位置するスコアをそれぞれのクラスタごとに得て、そのクラスタごとの中央または中央付近に位置するスコアをそれぞれのクラスタスコアとし、そのクラスタスコアによって、クラスタの順位を再構成するようにしてもよい。

## 【0025】

また、前記クラスタリング処理を複数の検索サービスによって得られた検索結果に対応して行うことを可能とする場合、前記クラスタの順位を再構成するためのクラスタスコアを求める処理は、複数の検索サービスそれぞれに対応して行うようにする。

## 【0026】

さらに、前記クラスタリング処理は、それぞれの文書のタイトルを検出し、そのタイトルに含まれる特徴的な用語を特徴要素として抽出し、抽出された特徴要素に基づいて行うようにしている。

## 【0027】

そして、前記クラスタ順位が再構成されたクラスタリング結果の出力の仕方は、クラスタスコアの高いクラスタ順に表示し、クラスタスコアが同じであるクラスタが存在する場合には、クラスタ内の文書数の多いクラスタを高順位として表



示する。

【0028】

このように本発明は、クラスタリング結果に対し、検索要求との適合性を示すスコアを用いて、それぞれのクラスタの順位を再構成して出力するようにしているので、ユーザは自分の欲しい情報を探しやすくなる。

【0029】

そして、それぞれのクラスタの順位を再構成する際に用いられるパラメータとして、クラスタごとのスコアの平均値を用いることで、それぞれのクラスタのキーワードに対する適合性を良好に反映した順位付けを行うことができる。

【0030】

また、スコアの平均値以外にもクラスタごとのスコアの最大値や中央値を用いることも可能である。最大値を用いることで、クラスタ順位を決めるために足し算や割り算を行う必要がなくなり、計算量を少なくすることができ、しかも、同じクラスタ内に、極端に低いスコアを持つ文書が少数あっても、その影響をあまり受けないようにすることができる。中央値も同様に、クラスタ順位を決めるための計算量を少なくすることができ、しかも、中央値の場合、同じクラスタ内に、極端に高いスコアや極端に低いスコアを持つ文書が少数あっても、その影響をあまり受けないようにすることができる。

【0031】

また、前記クラスタリング処理は、それぞれの文書のタイトルに含まれる特徴要素に基づいて行うことによって、少ない演算量で適切なクラスタリング処理を行うことができる。

【0032】

また、複数の検索サービスによって得られた検索結果に対応してクラスタリングを行う場合、クラスタスコアを求める処理を、複数の検索サービスそれぞれに対応して行うようにすれば、特定の検索サービスだけではなく、幾つもの検索サービスで検索された結果に対してもクラスタ順位の再構成処理を行うことができる。

## 【0033】

そして、クラスタ順位が再構成されたクラスタリング結果の出力の仕方としては、クラスタスコアの高いクラスタ順に表示し、さらに、クラスタスコアが同じであるクラスタが存在する場合には、クラスタ内の文書数の多いクラスタを高順位とし、それを一覧表形式で表示することにより、ユーザは自分の欲しい情報を探す際、その一覧表を上から順に見て行けばよいことから、効率よく情報をさがすことができる。

## 【0034】

## 【発明の実施の形態】

以下、本発明の実施の形態について説明する。なお、この実施の形態で説明する内容は、本発明の情報分類方法および情報分類装置についての説明であるとともに、本発明の情報分類処理プログラムを記録した記録媒体における情報分類処理プログラムの具体的な処理内容をも含むものである。

## 【0035】

図1は本発明を実現するための装置構成を示すもので、大きく分けると、検索サービス1、変換モジュール2、クラスタリングモジュール3、クラスタ順位再構成モジュール4とから構成され、変換モジュール2、クラスタリングモジュール3、クラスタ順位再構成モジュール4が本発明の情報分類装置に相当し、特に、クラスタ順位再構成モジュール4を設けたことに本発明の特徴がある。

## 【0036】

検索サービス1はインターネットなどで一般に広く使用されている汎用の検索サービスであり、ユーザからの検索要求としてのキーワードが入力されることにより、たとえばwebページなどから、入力されたキーワードに基づいた情報検索を行う。この検索サービス1で検索された検索結果はファイル形式で出力され、クラスタリングモジュール3に渡されるが、複数の検索サービスが存在する場合、それぞれの検索サービスによって出力されるデータ形式が異なるので、どのような検索サービスからのファイルであってもその内容を読めるような形式にファイルを変換するために変換モジュール2が設けられる。

## 【0037】

クラスタリングモジュール3は、検索サービス1により出力された検索結果（変換モジュール2による変換後のファイル内容）に対しクラスタリング処理を行うもので、この実施の形態では、それぞれの文書から文書のタイトルを抽出し、そのタイトルに含まれる特徴的な用語を特徴要素として抽出し、抽出された特徴要素に基づいてクラスタリング処理を行う。

## 【0038】

具体的には、それぞれの文書からタイトルとして抽出された部分を形態素解析し、形態素解析された結果から特徴的な用語を特徴要素として抽出する。その後、抽出された特徴要素とそれぞれの文書との関係を示す特徴テーブルを作成する。この特徴テーブルは、たとえば、抽出されたそれぞれの特徴要素が各文書のタイトルに幾つ含まれているかを、それぞれの特徴要素とそれぞれの文書と対応づけて示すもので、一例として、それぞれの文書のタイトルから、「概要」、「LP」、「仕様」、「デバイス」、「半導体」、「電子」というような特徴要素が抽出されたとすると、これらの特徴要素が、それぞれの文書のそれぞれのタイトルに、それぞれ何個含まれているかを示す内容となっている。

## 【0039】

このような特徴テーブルに基づいて、多数の文書を意味的に共通性のある複数のクラスタに分類する。つまり、それぞれの文書のそれぞれのタイトルに存在する特徴要素に基づいて、共通する特徴要素を持つ文書を1つのまとまりとし、そのまとまりを1つのクラスタとする。

## 【0040】

このクラスタリングモジュール3から、たとえば、図2のようなクラスタリング結果が出力されたとする。この図2は前述したように、クラスタリングされることによって得られた各クラスタの名称（上述の特徴要素に相当し、それをここではクラスタ名と呼んでいる）と、そのクラスタに属するそれぞれの文書のタイトルと、そのクラスタには幾つの文書が含まれているかを示す文書数、さらには、前述したように、それぞれのタイトルごとにスコアを示す数値などが一覧表形式で示されている。

## 【0041】

なお、このスコアは、前述したように、与えられたキーワードとそれぞれの文書との適合度を示す客観的な尺度として用いられ、ここでは、値が大きいほど、与えられたキーワードに適合する内容を有する文書であるとする。このスコアはキーワードとの適合度を表すものであるので、その単位としては、%や点数など検索サービスによって異なるがこの実施の形態では点数で表すものとする。

## 【0042】

そして、このクラスタリングモジュール3によってクラスタリングされた段階のクラスタリング結果は、図2に示されるように、クラスタの配置はそれぞれのクラスタに含まれる文書の数の多い順となっている。前述したように、この場合、上から順に、概要クラスタ、LPクラスタ、仕様クラスタ、デバイスクラスタ、半導体クラスタ、電子クラスタの順となっている。

## 【0043】

クラスタ順位再構成モジュール4は、クラスタリングモジュール3で出力されたクラスタリング結果に基づいて、それぞれのクラスタの表示順位を再構成するもので、その処理内容の詳細については後に説明する。

## 【0044】

このような構成において、本発明の情報分類処理について説明する。本発明が行う情報分類処理手順は概略的には、図3のフローチャートに示すように、まず、汎用の検索サービス1で検索された検索結果を取得し（ステップS1）、取得した検索結果に対しクラスタリング処理を施し（ステップS2）、そのクラスタリング結果を出力する（ステップS3）。そして、そのクラスタリング結果に対し、それぞれのクラスタ順位の再構成を行い（ステップS4）、再構成されたクラスタリング結果を出力する（ステップS5）。以下、具体例を参照しながら詳細に説明する。

## 【0045】

この実施の形態では、クラスタリングモジュール3が行うクラスタリング処理は、検索サービス1で検索された文書に対し、それぞれの文書のタイトルを抽出し、そのタイトルから特徴要素を抽出して、抽出された特徴要素とそれぞれの文

書との関係を示す特徴テーブルを作成して、その特徴テーブルの内容に基づいて、それぞれの文書を意味的に共通性のある複数のクラスタに分類する。また、この実施の形態では、ユーザが「半導体」というキーワードを検索要求として検索サービス 1 に入力し、それによって得られた多数の文書がクラスタリングモジュール 3 によってクラスタリングされ、そのクラスタリング結果が図 2 に示すような結果であったとする。

## 【0046】

このクラスタリングモジュール 3 からのクラスタリング結果は、クラスタ順位再構成モジュール 4 に入力され、以下に示すような処理がなされる。

## 【0047】

まず、図 2 で示されるクラスタリング結果における各クラスタ（概要クラスタ、LP クラスタ、仕様クラスタ、デバイスクラスタ、半導体クラスタ、電子クラスタ）において、それぞれのクラスタに含まれる文書対応に付されたスコアを利用して、そのスコアの値の平均を求める。この場合、それぞれのクラスタごとにスコアの値を足し算し、その足し算して得られた結果をそのクラスタに含まれる文書数で割る単純平均を求める。

## 【0048】

たとえば、概要クラスタで考えると、この図 2 に示す検索結果においては、そのクラスタ内のスコアの合計が 579 点あって、文書数が 16 個であるので、平均のスコアは約 36 点と求められる。また、「LP」クラスタで考えると、そのクラスタ内のスコアの合計が 450 点であって、文書数が 16 個であるので、平均のスコアは約 28 点と求められる。同様に、「仕様」クラスタは、そのクラスタ内のスコアの合計が 413 点であって、文書数が 14 個であるので、平均のスコアは約 29 点と求められ、「デバイス」クラスタは、そのクラスタ内のスコアの合計が 849 点であって、文書数が 9 個であるので、平均のスコアは約 94 点と求められ、「半導体」クラスタは、そのクラスタ内のスコアの合計が 757 点であって、文書数が 7 個であるので、平均のスコアは約 108 点と求められ、「電子」クラスタは、そのクラスタ内のスコアの合計が 349 点であって、文書数が 4 個であるので、平均のスコアは約 87 点と求められる。

## 【0049】

以上のようにして計算された平均のスコアを各クラスタのスコア（クラスタスコアと呼ぶ）とする。そして、このクラスタスコアの高い順にクラスタの順位を再構成する。

## 【0050】

すなわち、この場合、クラスタスコアの最も高いクラスタは、半導体クラスタの108点であり、第2位はデバイスクラスタの94点であり、第3位は電子クラスタの87点であり、以下、概要クラスタ（36点）、仕様クラスタ（29点）、LPクラスタ（28点）といった順序となる。

## 【0051】

このようにして、それぞれのクラスタごとにクラスタスコアを計算し、求められたクラスタスコアの高い順にクラスタ順位を再構成する。この再構成されたクラスタリング結果を一覧表形式で表したものが図4である。図4によれば、表の最上段に半導体クラスタが位置し、2番目にデバイスクラスタ、3番目に電子クラスタ、以下、概要クラスタ、仕様クラスタ、LPクラスタといった順序となる。この図4のクラスタリング結果によれば、ユーザの入力した「半導体」というキーワードに対し、そのキーワードに適合する文書が多く含まれるクラスタが上位に来ていることがわかる。

## 【0052】

この図4のクラスタリング結果と図2のクラスタリング結果を比較すると、図2のクラスタリング結果では、ユーザの入力した「半導体」というキーワードに対し、そのキーワードとは直接には関係しないような文書で構成される概要クラスタ、LPクラスタ、仕様クラスタといったクラスタが上位に位置し、キーワードに大きく関係するような文書が含まれると思われる半導体クラスタ、デバイスクラスタ、電子クラスタといったクラスタが下位に位置しているが、図4では、それが逆転し、キーワードに大きく関係するような文書が含まれると思われるクラスタが上位に位置するようになる。

## 【0053】

なお、クラスタスコアが同じ値となった場合には、クラスタ内に含まれる文書

数の多い方を上位とするなどの措置を講ずる。

【0054】

以上説明したように、単純にそれぞれのクラスタに含まれる文書数（1つのクラスタにまとめられた文書数）によって順位付けするのではなく、それぞれのクラスタごとにそのクラスタに含まれる文書に付されたスコアに基づいてクラスタの順位を決めることによって、キーワードに適合したクラスタ順位が得られる。

【0055】

なお、図4に示すようなクラスタリング結果がユーザに表示され、ユーザはこのようなクラスタリング結果の一覧表を見て、自分の欲しい情報の入っている文書のタイトル部分をクリックすれば、そのタイトルに対応する本文が表示されるというような表示処理がなされる。

【0056】

以上説明したように、この実施の形態では、ユーザの入力したキーワードによって検索された多数の文書に対し、これら多数の文書のタイトルに含まれる特徴要素に基づいてクラスタリング処理し、さらに、そのクラスタリング結果に対して、それぞれのクラスタごとにそのクラスタに含まれる文書のスコアの平均を求める。そして、その平均のスコアをクラスタスコアとし、それぞれのクラスタごとのクラスタスコアに基づいて、クラスタ順位の再構成を行う。つまり、クラスタスコアの大きい順にクラスタの並べ替えを行い、図4に示すようなクラスタリング結果として表示する。

【0057】

これによって、ユーザの欲しい情報の入っているクラスタが一覧表の上位に位置した状態で表示されているので、自分の欲しい情報を探しやすくなる。

【0058】

また、これまでの説明では、ある1つの汎用の検索サービスで検索された結果をクラスタリング処理する場合について説明したが、複数の検索サービスにより検索された結果をクラスタリング処理する場合にも適用できる。

【0059】

検索サービスはそれぞれに得意の分野があることも多く、たとえば、ある検索

サービスはスポーツ関係の情報を多数保有し、ある検索サービスは学術関係の情報を多数保有し、また、ある検索サービスは芸能関係の情報を多数保有しているというように、それぞれの得意の分野が存在する場合も多い。これらそれぞれの得意分野については豊富な情報を所有しており、ユーザの所望とする情報が適切に取り出される可能性が高い。したがって、情報検索を行う際は、検索サービスを使い分けることも普通に行われる。

#### 【0060】

このように、複数の検索サービスを用いてクラスタリング処理する場合には、それぞれの検索サービスにより検索された検索結果の内容、長さ、検索結果出力順序などがまちまちなので、それぞれの検索サービスからのファイルをクラスタリングモジュール3で処理可能な形式に変換する変換モジュール2を複数の検索サービスに対応して用意する。そして、さらに、そのクラスタリング結果におけるクラスタ順位再構成を行う場合には、それぞれのクラスタのクラスタスコアを求める処理をそれぞれの検索サービスに対応して行うようにする。

#### 【0061】

たとえば、本発明の要旨であるクラスタ順位再構成処理についていえば、検索サービスによって幾つかの対策を講じる必要がある。たとえば、スコアの分布の幅が非常に大きい場合（たとえば、スコアを表す数値が最大1000から最小は2など）は、対数を取って計算するなどの措置を講じたり、また、きわめてスコアの値が小さい文書（たとえば、殆どの文書が数百のスコアの値があるのに2や3の値しかない文書）はクラスタリング対象から外すといった措置を講じる。

#### 【0062】

このように、複数の検索サービスに対応できるようにすることで、ユーザは検索しようとする情報の分野に応じて検索サービスを使い分けることができ、それぞれの得意分野に応じた検索が可能となるばかりでなく、ある1つの検索サービスが混み合っているような場合には、他の検索サービスに切り換えて検索を行うというような柔軟な検索も可能となる。

#### 【0063】

なお、本発明は以上説明した実施の形態に限定されるものではなく、本発明の



要旨を逸脱しない範囲で種々変形実施可能となるものである。たとえば、これまで説明した実施の形態では、それぞれのクラスタのクラスタスコアは、そのクラスタに含まれる文書のスコアの単純平均を用いた例について説明したが、このクラスタスコアとしては、それぞれのクラスタ内に含まれる文書のなかで最大のスコアを有する文書のスコアを用いるようにしてもよく、また、それぞれのクラスタ内に含まれる文書に付されたスコアのなかで中央に位置する文書のスコアを用いるようにしてもよい。

#### 【0064】

このように、クラスタごとのスコアの最大値を用いることで、クラスタ順位を決めるために足し算したり割り算したりという計算を行う必要がなく、計算量を少なくすることができ、しかも、同じクラスタ内に、極端に低いスコアを持つ文書が少数あっても、その影響をあまり受けないようにすることができる。また、クラスタごとのスコアの中央値を用いる場合も、最大値を用いるのと同様、クラスタ順位を決めるための計算量を少なくすることができ、しかも、中央値の場合、同じクラスタ内に、極端に高いスコアや極端に低いスコアを持つ文書が少数あっても、その影響をあまり受けないようにすることができる。

#### 【0065】

また、前述の実施の形態では、クラスタリングを行うための情報（クラスタリング対象情報）として、検索されたそれぞれの文書のタイトルを用いた例について説明したが、これは、タイトルだけでなく、たとえば、URLアドレス（http://を取り除いた部分）、更新日時（単純な時間または最近1カ月以内の更新日時）、ファイルサイズ（webページ本文のバイトサイズなど）を用いてクラスタリングすることもできる。また、これらは、単独で用いてクラスタリングするようにしてもよく、幾つかを組み合わせてもよい。このように、クラスタリング対象情報を種々選ぶことによって、それぞれに応じた特色のあるクラスタリングが行える。そして、これらのどれを用いるかは、最初にメニューなどで選択項目を選ぶことで可能となる。また、選んだ項目が無い場合には、他の項目を代用する。たとえば、タイトルを選んだ場合、webページにタイトルが無い場合には、URLアドレスを代用する。

## 【0066】

また、以上説明した本発明の情報分類処理を行う処理プログラムは、フロッピーディスク、光ディスク、ハードディスクなどの記録媒体に記録しておくことができ、本発明はその記録媒体をも含むものである。また、ネットワークから処理プログラムを得るようにしてもよい。

## 【0067】

## 【発明の効果】

以上説明したように本発明によれば、検索された複数の文書をクラスタリング処理し、そのクラスタリング処理結果に対し、検索要求との適合性を示すスコアを用いて、それぞれのクラスタの順位を再構成して出力し、クラスタスコアの高い順に表形式で表示するようにしているので、ユーザは自分の欲しい情報を探す際、その一覧表を上から順に見て行けばよいことから、効率よく情報をさがすことができる。しかも、クラスタ順位再構成処理としては、もともと存在するスコアを用いて、それぞれのクラスタごとに平均値や最大値あるいは中央値などを求め、その結果に基づいて並べ替えするだけの処理であるので、複雑な演算処理が不要で短時間での処理が可能となる。これによって、汎用の検索サービスを利用してwebページの検索を行ったのちに、その検索結果の後処理として本発明を適用する場合にも、殆どリアルタイムで検索結果（クラスタ順位再構成後の最終的なクラスタリング結果）をユーザに提示することが可能となる。

## 【図面の簡単な説明】

## 【図1】

本発明の情報分類装置の実施の形態を説明する構成図である。

## 【図2】

ある検索サービスで検索された検索結果としての複数の文書をクラスタリングした結果の一例を示す図である。

## 【図3】

本発明の情報分類処理手順を概略的に説明するフローチャートである。

## 【図4】

図2で示されたクラスタリング結果をクラスタ順位再構成処理した結果を示す

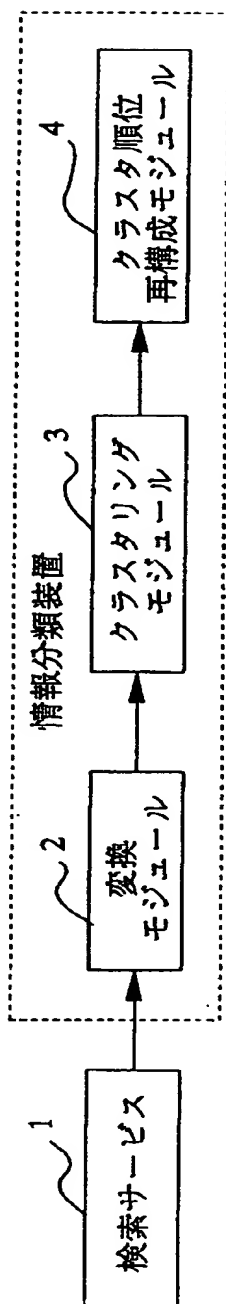
図である。

【符号の説明】

- 1 検索サービス
- 2 変換モジュール
- 3 クラスタリングモジュール
- 4 クラスタ順位再構成モジュール

【書類名】 図面

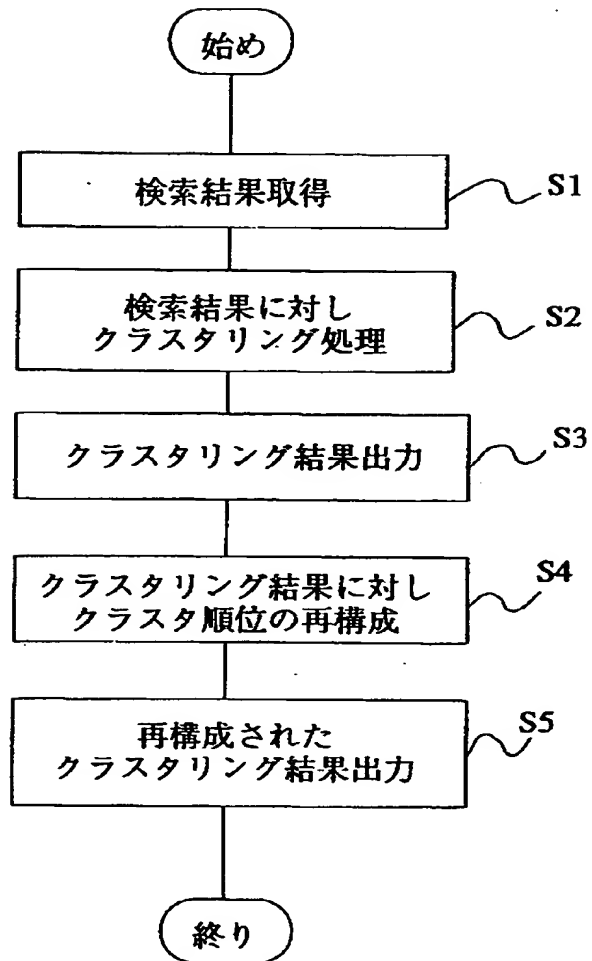
【図 1】



【図 2】

クラス名	文書数	スコア	文書のタイトル	検索結果
概要	16			133 東北エプソンの概要 49 LP-7000 仕様概要 38 LP-8000 仕様概要 38 LP-8400 仕様概要 34 LP-700 仕様概要 33 セイコーエプソン会社概要 30 LP-9200 仕様概要 28 LP-9200PS2 仕様概要 28 LP-830PS 仕様概要 27 LP-8300 仕様概要 26 LP-8600 仕様概要 24 LP-8200 仕様概要 23 LP-710 仕様概要 23 LP-800 仕様概要 23 LP-900 仕様概要 22 LP-500 仕様概要 (クラス内合計 579 平均 36)
LP	16			49 LP-7000 仕様概要 38 LP-8000 仕様概要 38 LP-8400 仕様概要 34 LP-700 仕様概要 30 LP-700S Spec Sheet & Option 30 LP-9200 仕様概要 28 LP-9200PS2 仕様概要 28 LP-830PS 仕様概要 27 LP-8300 仕様概要 26 LP-8600 仕様概要 24 LP-8200 仕様概要 23 LP-710 仕様概要 23 LP-800 仕様概要 23 LP-900 仕様概要 22 LP-500 仕様概要 7 エスパーレーザ LP-8300 (クラス内合計 450 平均 28)
仕様	14			49 LP-7000 仕様概要 38 LP-8000 仕様概要 38 LP-8400 仕様概要 34 LP-700 仕様概要 30 LP-9200 仕様概要 28 LP-9200PS2 仕様概要 28 LP-830PS 仕様概要 27 LP-8300 仕様概要 26 LP-8600 仕様概要 24 LP-8200 仕様概要 23 LP-710 仕様概要 23 LP-800 仕様概要 23 LP-900 仕様概要 22 LP-500 仕様概要 (クラス内合計 413 平均 29)
デバイス	9			117 デバイス 半導体 CARD-PC 115 デバイス 半導体 メモリ 111 電子デバイス ASIC 111 デバイス 半導体 PCカード製品 101 デバイス コラボレーション 90 電子デバイス 半導体 マイコン 82 EPSON 電子デバイス新製品 66 EPSON 電子デバイスお問い合わせ 56 デバイス 半導体 ASSP (クラス内合計 849 平均 94)
半導体	7			157 半導体事業部環境方針 117 デバイス 半導体 CARD-PC 115 デバイス 半導体 メモリ 111 電子デバイス ASIC 111 デバイス 半導体 PCカード製品 90 電子デバイス 半導体 マイコン 56 デバイス 半導体 ASSP (クラス内合計 757 平均 108)
電子	4			111 電子デバイス ASIC 90 電子デバイス 半導体 マイコン 82 EPSON 電子デバイス新製品 66 EPSON 電子デバイスお問い合わせ (クラス内合計 349 平均 87)

【図 3】



【図 4】

クラス名	文書数	スコア	文書のタイトル	検索結果
半導体	7			157 半導体事業部環境方針 117 デバイス 半導体 CARD-PC 115 デバイス 半導体 メモリ 111 電子デバイス ASIC 111 デバイス 半導体 PCカード製品 90 電子デバイス 半導体 マイコン 56 デバイス 半導体 ASSP (クラス内合計 757 平均 108)
デバイス	9			117 デバイス 半導体 CARD-PC 115 デバイス 半導体 メモリ 111 電子デバイス ASIC 111 デバイス 半導体 PCカード製品 101 デバイス コラボレーション 90 電子デバイス 半導体 マイコン 82 EPSON 電子デバイス新製品 66 EPSON 電子デバイスお問い合わせ 56 デバイス 半導体 ASSP (クラス内合計 849 平均 94)
電子	4			111 電子デバイス ASIC 90 電子デバイス 半導体 マイコン 82 EPSON 電子デバイス新製品 66 EPSON 電子デバイスお問い合わせ (クラス内合計 349 平均 87)
概要	16			133 東北エプソンの概要 49 LP-7000 仕様概要 38 LP-8000 仕様概要 38 LP-8400 仕様概要 34 LP-700 仕様概要 33 セイコーエプソン会社概要 30 LP-9200 仕様概要 28 LP-9200PS2 仕様概要 28 LP-830PS 仕様概要 27 LP-8300 仕様概要 26 LP-8600 仕様概要 24 LP-8200 仕様概要 23 LP-710 仕様概要 23 LP-800 仕様概要 23 LP-900 仕様概要 22 LP-500 仕様概要 (クラス内合計 579 平均 36)
仕様	14			49 LP-7000 仕様概要 38 LP-8000 仕様概要 38 LP-8400 仕様概要 34 LP-700 仕様概要 30 LP-9200 仕様概要 28 LP-9200PS2 仕様概要 28 LP-830PS 仕様概要 27 LP-8300 仕様概要 26 LP-8600 仕様概要 24 LP-8200 仕様概要 23 LP-710 仕様概要 23 LP-800 仕様概要 23 LP-900 仕様概要 22 LP-500 仕様概要 (クラス内合計 413 平均 29)
LP	16			49 LP-7000 仕様概要 38 LP-8000 仕様概要 38 LP-8400 仕様概要 34 LP-700 仕様概要 30 LP-700S Spec Sheet & Option 30 LP-9200 仕様概要 28 LP-9200PS2 仕様概要 28 LP-830PS 仕様概要 27 LP-8300 仕様概要 26 LP-8600 仕様概要 24 LP-8200 仕様概要 23 LP-710 仕様概要 23 LP-800 仕様概要 23 LP-900 仕様概要 22 LP-500 仕様概要 7 エスパーレーザー LP-8300 (クラス内合計 450 平均 28)

【要約】

【課題】 汎用の検索サービスによって検索された多数の文書をクラスタリングし、それを一覧表形式で表示する際、そのクラスタリングによって得られたクラスターの並びが、ユーザの与えたキーワードにあまり関係のない文書を含むクラスターが上位に位置してしまう場合がある。

【解決手段】 汎用の検索サービス1で検索された複数の検索結果をクラスタリングモジュール3が取得して、検索結果をクラスタリング処理し、さらに、そのクラスタリング結果に対し、クラスタ順位再構成ジュール4がそれぞれのクラスタに含まれるそれぞれの文書対応に付された前記検索要求との適合性を示すスコアを用い、クラスタごとにスコアの平均を求めるなどしてクラスタごとのスコア（クラスタスコア）を求め、それによってクラスタの順位を再構成し、クラスタスコアの高い順に表形式で出力する。

【選択図】

図1



出 願 人 履 歴 情 報

識別番号 [000002369]

1. 変更年月日	1990年 8月20日
[変更理由]	新規登録
住 所	東京都新宿区西新宿2丁目4番1号
氏 名	セイコーエプソン株式会社

**THIS PAGE BLANK (USPTO)**